

# Incorporating Generative AI in Model Governance Programs

Patrick Hall

June 4, 2024

## Abstract

## 1 Introduction

The National Institute of Standards and Technology Artificial Intelligence (AI) Risk Management Framework (RMF).[\[25\]](#)

## 2 Generative AI Incidents

## 3 Generative AI Governance

## 4 Generative AI Inventories

## 5 Generative AI Risk Tiers

## 6 Generative AI Risk Measurement

## 7 Generative AI Risk Management

## Conclusion

## Acknowledgments

Thank you to Bernie Siskin and Nick Schmidt of BLDS and Eric Sublett of Relman Colfax for formative discussions relating to GAI risk tiering.

## Abbreviations

- AI: Artificial Intelligence
- AI RMF: Artificial Intelligence Risk Management Framework
- GAI: Generative AI
- LLM: Large Language Model
- RMF: Risk Management Framework

- [1] Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabisa. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. *arXiv preprint arXiv:2308.16884*, 2023.
- [2] Rishi Bommasani, Percy Liang, and Tony Lee. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*, 1525(1):140–146, 2023.
- [3] Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.
- [4] Adrian de Wynter, Xun Wang, Alex Sokolov, Qilong Gu, and Si-Qing Chen. An evaluation on large language model outputs: Discourse and memorization. *Natural Language Processing Journal*, 4:100024, 2023.
- [5] Innovation Department for Science and Technology. International scientific report on the safety of advanced ai. *gov.uk*, 2024. <https://www.gov.uk/government/publications/international-scientific-report-on-the-safety-of-advanced-ai>.
- [6] Jeremy Dohmann. Blazingly fast llm evaluation for in-context learning. <https://www.databricks.com/blog/llm-evaluation-for-icl>. Last accessed: May 24, 2024.
- [7] Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do membership inference attacks work on large language models? *arXiv:2402.07841*, 2024.
- [8] Esin Durmus, Karina Nyugen, Thomas I Liao, Nicholas Schiefer, Amanda Askill, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, et al. Towards measuring the representation of subjective global opinions in language models. *arXiv preprint arXiv:2306.16388*, 2023.
- [9] Hugging Face. Evaluation. <https://huggingface.co/docs/evaluate/index>. Last accessed: May 24, 2024.
- [10] Shangbin Feng, Chan Young Park, Yuhang Liu, and Yulia Tsvetkov. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair nlp models. *arXiv preprint arXiv:2305.08283*, 2023.
- [11] Jack FitzGerald, Christopher Hench, Charith Peris, Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron Nash, Liam Urbach, Vishesh Kakarala, Richa Singh, et al. Massive: A 1m-example multilingual natural language understanding dataset with 51 typologically-diverse languages. *arXiv preprint arXiv:2204.08582*, 2022.
- [12] AI Verify Foundation. Cataloguing llm evaluations. <https://aiverifyfoundation.sg/>, 2023. See also: <https://github.com/aiverify-foundation/LLM-Evals-Catalogue>.
- [13] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 12 2023.
- [14] Patrick Hall and Daniel Atherton. Awesome machine learning interpretability, 2024. <https://github.com/jphall663/awesome-machine-learning-interpretability>.
- [15] Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37, 2022.

- [16] Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. Catastrophic jailbreak of open-source llms via exploiting generation. In *The Twelfth International Conference on Learning Representations*, 2023.
- [17] Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [18] ISO. Information technology — artificial intelligence — management system. *ISO/IEC 42001:2023*, 2023. <https://www.iso.org/obp/ui/en/#iso:std:iso-iec:42001:ed-1:v1:en>.
- [19] Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024.
- [20] Gary McGraw, Harold Figueroa, Katie McMahon, and Richie Bonett. An architectural risk analysis of large language models: Applied machine learning security. *Berryville Inst. Mach. Learn.(BIML), Berryville, VA, USA, Tech. Rep*, 2024.
- [21] Gary McGraw, Harold Figueroa, Victor Shepardson, and Richie Bonett. An architectural risk analysis of machine learning systems: Toward more secure machine learning. *Berryville Institute of Machine Learning, Clarke County, VA. Accessed on: Mar, 23, 2020*.
- [22] Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*, 2023.
- [23] Microsoft. Microsoft responsible ai standard, v2. <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RE5cmFl>.
- [24] NIST. Guide for conducting risk assessments. *SP800-03R1*, pages i–L2, 2012.
- [25] NIST. Artificial Intelligence Risk Management Framework (AI RMF 1.0). *nist.gov*, 2023.
- [26] NIST. Nist ai rmf playbook, 2023. [https://airc.nist.gov/AI\\_RMF\\_Knowledge\\_Base/Playbook](https://airc.nist.gov/AI_RMF_Knowledge_Base/Playbook).
- [27] NIST. Ai 600-1, generative ai risk profile (draft). *nist.gov*, 2024. <https://airc.nist.gov/docs/NIST.AI.600-1.GenAI-Profile.ipd.pdf>.
- [28] Office of the Comptroller of the Currency. Model risk management. *OCC Handbook*, 2021. <https://www.occ.gov/publications-and-resources/publications/comptrollers-handbook/files/model-risk-management/index-model-risk-management.html>.
- [29] Julien Piet, Chawin Sitawarin, Vivian Fang, Norman Mu, and David Wagner. Mark my words: Analyzing and evaluating language model watermarks. *arXiv preprint arXiv:2312.00273*, 2023.
- [30] Jérôme Rutinowski, Sven Franke, Jan Endendyk, Ina Dormuth, Moritz Roidl, Markus Pauly, et al. The self-perception and political biases of chatgpt. *Human Behavior and Emerging Technologies*, 2024, 2023.
- [31] Elvis Saravia. Prompt Engineering Guide. <https://github.com/dair-ai/Prompt-Engineering-Guide>, 12 2022.
- [32] Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. *arXiv preprint arXiv:2308.03825*, 2023.
- [33] Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*, 2023.

- [34] Chawin Sitawarin and Charlie Cheng-Jie Ji. Llm security & privacy. <https://github.com/chawins>, 2024.
- [35] Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. "i'm sorry to hear that": Finding new biases in language models with a holistic descriptor dataset. *arXiv preprint arXiv:2205.09209*, 2022.
- [36] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022.
- [37] Robin Staab, Mark Vero, Mislav Balunović, and Martin Vechev. Beyond memorization: Violating privacy via inference with large language models. *arXiv preprint arXiv:2310.07298*, 2023.
- [38] Victor Storchon, Ravin Kumar, Rumman Chowdhury, Seraphina Goldfarb-Tarrant, and Sven Cattell. Generative ai red teaming challenge: Transparency report. <https://www.humane-intelligence.org/>, 2024.
- [39] Bertie Vidgen, Adarsh Agrawal, Ahmed M Ahmed, Victor Akinwande, Namir Al-Nuaimi, Najla Alfaraj, Elie Alhajjar, Lora Aroyo, Trupti Bavalatti, Borhane Blili-Hamelin, et al. Introducing v0. 5 of the ai safety benchmark from mlcommons. *arXiv preprint arXiv:2404.12241*, 2024.
- [40] Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [41] Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. Flask: Fine-grained language model evaluation based on alignment skill sets. *arXiv preprint arXiv:2307.10928*, 2023.
- [42] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.

# Appendix A: Example Generative AI–Trustworthy Characteristic Crosswalk

## A.1: Trustworthy Characteristic to Generative AI Risk Crosswalk

Table A.1: Trustworthy Characteristic to Generative AI Risk Crosswalk.

| Accountable and Transparent  | Explainable and Interpretable                                   | Fair with Harmful Bias Managed   | Privacy Enhanced   |
|--|---|--|--|
| Data Privacy<br>Environmental<br>Human-AI Configuration<br>Information Integrity<br>Intellectual Property<br>Value Chain and Component Integration | Human-AI Configuration<br>Value Chain and Component Integration | Confabulation<br>Environmental<br>Human-AI Configuration<br>Intellectual Property<br>Obscene, Degrading, and/or Abusive Content<br>Toxicity, Bias, and Homogenization<br>Value Chain and Component Integration | Data Privacy<br>Human-AI Configuration<br>Information Security<br>Intellectual Property<br>Value Chain and Component Integration |

  

| Safe   | Secure and Resilient  | Valid and Reliable  |
|--|---|---|
| CBRN Information<br>Confabulation<br>Dangerous or Violent Recommendations<br>Data Privacy<br>Environmental<br>Human-AI Configuration<br>Information Integrity<br>Information Security<br>Obscene, Degrading, and/or Abusive Content<br>Value Chain and Component Integration | Dangerous or Violent Recommendations<br>Data Privacy<br>Human-AI Configuration<br>Information Security<br>Value Chain and Component Integration | Confabulation<br>Human-AI Configuration<br>Information Integrity<br>Information Security<br>Toxicity, Bias, and Homogenization<br>Value Chain and Component Integration |

## A.2: Generative AI Risk to Trustworthy Characteristic Crosswalk

Table A.2: Generative AI Risk to Trustworthy Characteristic Crosswalk.

| CBRN Information | Confabulation  | Dangerous or Violent Recommendations | Data Privacy  |
|------------------|--|--------------------------------------|---|
| Safe             | Fair with Harmful Bias Managed<br>Safe<br>Valid and Reliable | Safe<br>Secure and Resilient         | Accountable and Transparent<br>Privacy Enhanced<br>Safe<br>Secure and Resilient |

  

| Environmental   | Human-AI Configuration   | Information Integrity                                     | Information Security   |
|---|--|---|--|
| Accountable and Transparent<br>Fair with Harmful Bias Managed<br>Safe | Accountable and Transparent<br>Explainable and Interpretable<br>Fair with Harmful Bias Managed<br>Privacy Enhanced<br>Safe<br>Secure and Resilient<br>Valid and Reliable | Accountable and Transparent<br>Safe<br>Valid and Reliable | Privacy Enhanced<br>Safe<br>Secure and Resilient<br>Valid and Reliable |

  

| Intellectual Property   | Obscene, Degrading, and/or Abusive Content | Toxicity, Bias, and Homogenization                   | Value Chain and Component Integration  |
|---|--|--|--|
| Accountable and Transparent<br>Fair with Harmful Bias Managed<br>Privacy Enhanced | Fair with Harmful Bias Managed<br>Safe     | Fair with Harmful Bias Managed<br>Valid and Reliable | Accountable and Transparent<br>Explainable and Interpretable<br>Fair with Harmful Bias Managed<br>Privacy Enhanced<br>Safe<br>Secure and Resilient<br>Valid and Reliable |

# Appendix B: Example Risk-tiering Materials for Generative AI

## B.1: Example Adverse Impacts

Table B.1: Example adverse impacts, adapted from NIST 800-30r1 Table H-2 [24].

| Level                       | Description   |
|-----------------------------|---|
| Harm to Operations          | <ul style="list-style-type: none"><li>• Inability to perform current missions/business functions.<ul style="list-style-type: none"><li>– In a sufficiently timely manner.</li><li>– With sufficient confidence and/or correctness.</li><li>– Within planned resource constraints.</li></ul></li><li>• Inability, or limited ability, to perform missions/business functions in the future.<ul style="list-style-type: none"><li>– Inability to restore missions/business functions.</li><li>– In a sufficiently timely manner.</li><li>– With sufficient confidence and/or correctness.</li><li>– Within planned resource constraints.</li></ul></li><li>• Harms (e.g., financial costs, sanctions) due to noncompliance.<ul style="list-style-type: none"><li>– With applicable laws or regulations.</li><li>– With contractual requirements or other requirements in other binding agreements (e.g., liability).</li></ul></li><li>• Direct financial costs.</li><li>• Reputational harms.<ul style="list-style-type: none"><li>– Damage to trust relationships.</li><li>– Damage to image or reputation (and hence future or potential trust relationships).</li></ul></li></ul> |
| Harm to Assets              | <ul style="list-style-type: none"><li>• Damage to or loss of physical facilities.</li><li>• Damage to or loss of information systems or networks.</li><li>• Damage to or loss of information technology or equipment.</li><li>• Damage to or loss of component parts or supplies.</li><li>• Damage to or loss of information assets.</li><li>• Loss of intellectual property.</li></ul>   |
| Harm to Individuals         | <ul style="list-style-type: none"><li>• Injury or loss of life.</li><li>• Physical or psychological mistreatment.</li><li>• Identity theft.</li><li>• Loss of personally identifiable information.</li><li>• Damage to image or reputation.</li><li>• Infringement of intellectual property rights.</li><li>• Financial harm or loss of income.</li></ul>   |
| Harm to Other Organizations | <ul style="list-style-type: none"><li>• Harms (e.g., financial costs, sanctions) due to noncompliance.<ul style="list-style-type: none"><li>– With applicable laws or regulations.</li><li>– With contractual requirements or other requirements in other binding agreements (e.g., liability).</li></ul></li><li>• Direct financial costs.</li><li>• Reputational harms.<ul style="list-style-type: none"><li>– Damage to trust relationships.</li><li>– Damage to image or reputation (and hence future or potential trust relationships).</li></ul></li></ul>  |
| Harm to the Nation          | <ul style="list-style-type: none"><li>• Damage to or incapacitation of critical infrastructure.</li><li>• Loss of government continuity of operations.</li><li>• Reputational harms.<ul style="list-style-type: none"><li>– Damage to trust relationships with other governments or with nongovernmental entities.</li><li>– Damage to national reputation (and hence future or potential trust relationships).</li></ul></li><li>• Damage to current or future ability to achieve national objectives.<ul style="list-style-type: none"><li>– Harm to national security.</li></ul></li><li>• Large-scale economic or workforce displacement.</li></ul>   |

## B.2: Example Impact Descriptions

Table B.2: Example Impact level descriptions, adapted from NIST SP800-30r1 Appendix H, Table H-3 [24].

| Qualitative Values | Semi-Quantitative Values |    | Description   |
|--------------------|--------------------------|----|---|
| Very High          | 96-100                   | 10 | An incident could be expected to have multiple severe or catastrophic adverse effects on organizational operations, organizational assets, individuals, other organizations, or the Nation.   |
| High               | 80-95                    | 8  | An incident could be expected to have a severe or catastrophic adverse effect on organizational operations, organizational assets, individuals, other organizations, or the Nation. A severe or catastrophic adverse effect means that, for example, the incident might: (i) cause a severe degradation in or loss of mission capability to an extent and duration that the organization is not able to perform one or more of its primary functions; (ii) result in major damage to organizational assets; (iii) result in major financial loss; or (iv) result in severe or catastrophic harm to individuals involving loss of life or serious life-threatening injuries.                       |
| Moderate           | 21-79                    | 5  | An incident could be expected to have a serious adverse effect on organizational operations, organizational assets, individuals other organizations, or the Nation. A serious adverse effect means that, for example, the incident might: (i) cause a significant degradation in mission capability to an extent and duration that the organization is able to perform its primary functions, but the effectiveness of the functions is significantly reduced; (ii) result in significant damage to organizational assets; (iii) result in significant financial loss; or (iv) result in significant harm to individuals that does not involve loss of life or serious life-threatening injuries. |
| Low                | 5-20                     | 2  | An incident could be expected to have a limited adverse effect on organizational operations, organizational assets, individuals other organizations, or the Nation. A limited adverse effect means that, for example, the incident might: (i) cause a degradation in mission capability to an extent and duration that the organization is able to perform its primary functions, but the effectiveness of the functions is noticeably reduced; (ii) result in minor damage to organizational assets; (iii) result in minor financial loss; or (iv) result in minor harm to individuals.  |
| Very Low           | 0-4                      | 0  | An incident could be expected to have a negligible adverse effect on organizational operations, organizational assets, individuals other organizations, or the Nation.  |



### B.3: Example Likelihood Descriptions

Table B.3: Example likelihood levels, adapted from NIST SP800-30r1 Appendix G, Table G-3 [24].

| Qualitative Values | Semi-Quantitative Values |    | Description   |
|--------------------|--------------------------|----|---|
| Very High          | 96-100                   | 10 | An incident is almost certain to occur; or occurs more than 100 times a year.                         |
| High               | 80-95                    | 8  | An incident is highly likely to occur; or occurs between 10-100 times a year.                         |
| Moderate           | 21-79                    | 5  | An incident is somewhat likely to occur; or occurs between 1-10 times a year.                         |
| Low                | 5-20                     | 2  | An incident is unlikely to occur; or occurs less than once a year, but more than once every 10 years. |
| Very Low           | 0-4                      | 0  | An incident is highly unlikely to occur; or occurs less than once every 10 years.                     |

### B.4: Example Risk Tiers

Table B.4: Example risk assessment matrix with 5 impact levels, 5 likelihood levels, and 5 risk tiers, adapted from NIST SP800-30r1 Appendix I, Table I-2 [24].

| Likelihood       | Level of Impact   |                   |                   |                   |                    |
|------------------|-------------------|-------------------|-------------------|-------------------|--------------------|
|                  | Very Low          | Low               | Moderate          | High              | Very High          |
| <b>Very High</b> | Very Low (Tier 5) | Low (Tier 4)      | Moderate (Tier 3) | High (Tier 2)     | Very High (Tier 1) |
| <b>High</b>      | Very Low (Tier 5) | Low (Tier 4)      | Moderate (Tier 3) | High (Tier 2)     | Very High (Tier 1) |
| <b>Moderate</b>  | Very Low (Tier 5) | Low (Tier 4)      | Moderate (Tier 3) | Moderate (Tier 3) | High (Tier 2)      |
| <b>Low</b>       | Very Low (Tier 5) | Low (Tier 4)      | Low (Tier 4)      | Low (Tier 4)      | Moderate (Tier 3)  |
| <b>Very Low</b>  | Very Low (Tier 5) | Very Low (Tier 5) | Very Low (Tier 5) | Low (Tier 4)      | Low (Tier 4)       |

## B.5: Example Risk Descriptions

Table B.5: Example risk descriptions, adapted from NIST SP800-30r1 Appendix I, Table I-3 [24] .

| Qualitative Values | Semi-Quantitative Values |    | Description   |
|--------------------|--------------------------|----|---|
| Very High          | 96-100                   | 10 | Very high risk means that an incident could be expected to have multiple severe or catastrophic adverse effects on organizational operations, organizational assets, individuals, other organizations, or the Nation. |
| High               | 80-95                    | 8  | High risk means that an incident could be expected to have a severe or catastrophic adverse effect on organizational operations, organizational assets, individuals, other organizations, or the Nation.              |
| Moderate           | 21-79                    | 5  | Moderate risk means that an incident could be expected to have a serious adverse effect on organizational operations, organizational assets, individuals, other organizations, or the Nation.                         |
| Low                | 5-20                     | 2  | Low risk means that an incident could be expected to have a limited adverse effect on organizational operations, organizational assets, individuals, other organizations, or the Nation.                              |
| Very Low           | 0-4                      | 0  | Very low risk means that an incident could be expected to have a negligible adverse effect on organizational operations, organizational assets, individuals, other organizations, or the Nation.                      |

## B.6: Practical Risk-tiering Questions

**B.6.1: Confabulation:** How likely are system outputs to contain errors? What are the impacts if errors occur?

**B.6.2: Dangerous and Violent Recommendations:** How likely is the system to give dangerous or violent recommendations? What are the impacts if it does?

**B.6.3: Data Privacy:** How likely is someone to enter sensitive data into the system? What are the impacts if this occurs? Are standard data privacy controls applied to the system to mitigate potential adverse impacts?

**B.6.4: Human-AI Configuration:** How likely is someone to use the system incorrectly or abuse it? How likely is use for decision-making? What are the impacts of incorrect use or abuse? What are the impacts of invalid or unreliable decision-making?

**B.6.5: Information Integrity:** How likely is the system to generate deepfakes or mis or disinformation? At what scale? Are content provenance mechanisms applied to system outputs? What are the impacts of generating deepfakes or mis or disinformation? Without controls for content provenance?

**B.6.6: Information Security:** How likely are system resources to be breached or exfiltrated? How likely is the system to be used in the generation of phishing or malware content? What are the impacts in these cases? Are standard information security controls applied to the system to mitigate potential adverse impacts?

**B.6.7: Intellectual Property:** How likely are system outputs to contain other entities' intellectual property? What are the impacts if this occurs?

**B.6.8: Toxicity, Bias, and Homogenization:** How likely are system outputs to be biased, toxic, homogenizing or otherwise obscene? How likely are system outputs to be used as subsequent training inputs? What are the impacts of these scenarios? Are standard nondiscrimination controls applied to mitigate potential adverse impacts? Is the application accessible to all user groups? What are the impacts if the system is not accessible to all user groups?

**B.6.9: Value Chain and Component Integration:** Are contracts relating to the system reviewed for legal risks? Are standard acquisition/procurement controls applied to mitigate potential adverse impacts? Do vendors provide incident response with guaranteed response times? What are the impacts if these conditions are not met?

# Appendix C: List of Selected Model Testing Suites

[12]

## C.1: Selected Model Testing Suites Organized by Trustworthy Characteristic

Table C.1: Selected model testing suites organized by trustworthy characteristic.

|   |
|---|
| <b>Accountable and Transparent</b>  |
| An Evaluation on Large Language Model Outputs:<br>Discourse and Memorization (see Appendix B)[4]<br>Big-bench: Truthfulness [36]<br>DecodingTrust: Machine Ethics [40]<br>Evaluation Harness: ETHICS [13]<br>HELM: Copyright [2]<br>Mark My Words [29]  |
| <b>Fair with Harmful Bias Managed</b>   |
| BELEBELE [1]<br>Big-bench: Low-resource language, Non-English, Translation<br>Big-bench: Social bias, Racial bias, Gender bias, Religious bias<br>Big-bench: Toxicity<br>DecodingTrust: Fairness<br>DecodingTrust: Stereotype Bias<br>DecodingTrust: Toxicity<br>C-Eval (Chinese evaluation suite) [17]<br>Evaluation Harness: CrowS-Pairs<br>Evaluation Harness: ToxiGen<br>Finding New Biases in Language Models with a Holistic Descriptor Dataset [35]<br>From Pretraining Data to Language Models to Downstream Tasks:<br>Tracking the Trails of Political Biases Leading to Unfair NLP Models [10]<br>HELM: Bias<br>HELM: Toxicity<br>MT-bench [42]<br>The Self-Perception and Political Biases of ChatGPT [30]<br>Towards Measuring the Representation of<br>Subjective Global Opinions in Language Models [8] |
| <b>Privacy Enhanced</b>   |
| HELM: Copyright<br>llmprivacy [37]<br>mimir [7]   |
| <b>Safe</b>   |
| Big-bench: Convince Me<br>Big-bench: Truthfulness<br>HELM: Reiteration, Wedging<br>Mark My Words<br>MLCommons [39]<br>The WMDP Benchmark [19]   |

Table C.1: Selected model testing suites organized by trustworthy characteristic (continued).

|  |
|--|
| <b>Secure and Resilient</b>  |
| Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation [16]  |
| DecodingTrust: Adversarial Robustness,<br>Robustness Against Adversarial Demonstrations  |
| detect-pretrain-code [33]  |
| In-The-Wild Jailbreak Prompts on LLMs [32]   |
| JailbreakingLLMs [3]   |
| llmprivacy   |
| mimir  |
| TAP: A Query-Efficient Method for Jailbreaking Black-Box LLMs [22]   |
| <b>Valid and Reliable</b>  |
| Big-bench: Algorithms, Logical reasoning, Implicit reasoning, Mathematics, Arithmetic, Algebra, Mathematical proof,<br>Fallacy, Negation, Computer code, Probabilistic reasoning, Social reasoning, Analogical reasoning, Multi-step,<br>Understanding the World |
| Big-bench: Analytic entailment, Formal fallacies and syllogisms with negation, Entailed polarity   |
| Big-bench: Context Free Question Answering   |
| Big-bench: Contextual question answering, Reading comprehension, Question generation   |
| Big-bench: Morphology, Grammar, Syntax   |
| Big-bench: Out-of-Distribution   |
| Big-bench: Paraphrase  |
| Big-bench: Sufficient information  |
| Big-bench: Summarization   |
| DecodingTrust: Out-of-Distribution Robustness, Adversarial Robustness, Robustness Against Adversarial Demonstrations   |
| Eval Gauntlet: Reading comprehension [6]   |
| Eval Gauntlet: Commonsense reasoning, Symbolic problem solving, Programming  |
| Eval Gauntlet: Language Understanding  |
| Eval Gauntlet: World Knowledge   |
| Evaluation Harness: BLiMP  |
| Evaluation Harness: CoQA, ARC  |
| Evaluation Harness: GLUE   |
| Evaluation Harness: HellaSwag, OpenBookQA, TruthfulQA  |
| Evaluation Harness: MuTual   |
| Evaluation Harness: PIQA, PROST, MC-TACO, MathQA, LogiQA, DROP   |
| FLASK: Logical correctness, Logical robustness, Logical efficiency, Comprehension, Completeness [41]   |
| FLASK: Readability, Conciseness, Insightfulness  |
| HELM: Knowledge  |
| HELM: Language   |
| HELM: Text classification  |
| HELM: Question answering   |
| HELM: Reasoning  |
| HELM: Robustness to contrast sets  |
| HELM: Summarization  |
| Hugging Face: Fill-mask, Text generation [9]   |
| Hugging Face: Question answering   |
| Hugging Face: Summarization  |
| Hugging Face: Text classification, Token classification, Zero-shot classification  |
| MASSIVE [11]   |
| MT-bench   |

## C.2: Selected Model Testing Suites Organized by Generative AI Risk

Table C.2: Selected model testing suites by organized generative AI risk.

|  |
|--|
| <b>CBRN Information</b>  |
| Big-bench: Convince Me   |
| Big-bench: Truthfulness  |
| HELM: Reiteration, Wedging   |
| MLCommons  |
| The WMDP Benchmark   |
| <b>Confabulation</b>   |
| BELEBELE   |
| Big-bench: Algorithms, Logical reasoning, Implicit reasoning, Mathematics, Arithmetic, Algebra, Mathematical proof, Fallacy, Negation, Computer code, Probabilistic reasoning, Social reasoning, Analogical reasoning, Multi-step, Understanding the World |
| Big-bench: Analytic entailment, Formal fallacies and syllogisms with negation, Entailed polarity   |
| Big-bench: Context Free Question Answering   |
| Big-bench: Contextual question answering, Reading comprehension, Question generation   |
| Big-bench: Convince Me   |
| Big-bench: Low-resource language, Non-English, Translation   |
| Big-bench: Morphology, Grammar, Syntax   |
| Big-bench: Out-of-Distribution   |
| Big-bench: Paraphrase  |
| Big-bench: Sufficient information  |
| Big-bench: Summarization   |
| Big-bench: Truthfulness  |
| C-Eval (Chinese evaluation suite)  |
| DecodingTrust: Out-of-Distribution Robustness, Adversarial Robustness, Robustness Against Adversarial Demonstrations   |
| Eval Gauntlet Reading comprehension  |
| Eval Gauntlet: Commonsense reasoning, Symbolic problem solving, Programming  |
| Eval Gauntlet: Language Understanding  |
| Eval Gauntlet: World Knowledge   |
| Evaluation Harness: BLiMP  |
| Evaluation Harness: CoQA, ARC  |
| Evaluation Harness: GLUE   |
| Evaluation Harness: HellaSwag, OpenBookQA, TruthfulQA  |
| Evaluation Harness: MuTual   |
| Evaluation Harness: PIQA, PROST, MC-TACO, MathQA, LogiQA, DROP   |
| FLASK: Logical correctness, Logical robustness, Logical efficiency, Comprehension, Completeness  |
| FLASK: Readability, Conciseness, Insightfulness  |
| Finding New Biases in Language Models with a Holistic Descriptor Dataset   |
| HELM: Knowledge  |
| HELM: Language   |
| HELM: Language (Twitter AAE)   |
| HELM: Question answering   |
| HELM: Reasoning  |
| HELM: Reiteration, Wedging   |
| HELM: Robustness to contrast sets  |
| HELM: Summarization  |
| HELM: Text classification  |
| Hugging Face: Fill-mask, Text generation   |
| Hugging Face: Question answering   |
| Hugging Face: Summarization  |
| Hugging Face: Text classification, Token classification, Zero-shot classification  |
| MASSIVE  |
| MLCommons  |
| MT-bench   |

Table C.2: Selected model testing suites by organized generative AI risk (continued).

|  |
|--|
| <b>Dangerous or Violent Recommendations</b>  |
| Big-bench: Convince Me   |
| Big-bench: Toxicity  |
| DecodingTrust: Adversarial Robustness, Robustness Against Adversarial Demonstrations                                 |
| DecodingTrust: Machine Ethics  |
| DecodingTrust: Toxicity  |
| Evaluation Harness: ToxiGen  |
| HELM: Reiteration, Wedging   |
| HELM: Toxicity   |
| MLCommons  |
| <b>Data Privacy</b>  |
| An Evaluation on Large Language Model Outputs: Discourse and Memorization (with human scoring, see Appendix B)       |
| Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation   |
| DecodingTrust: Machine Ethics  |
| Evaluation Harness: ETHICS   |
| HELM: Copyright  |
| In-The-Wild Jailbreak Prompts on LLMs  |
| JailbreakingLLMs   |
| MLCommons  |
| Mark My Words  |
| TAP: A Query-Efficient Method for Jailbreaking Black-Box LLMs  |
| detect-pretrain-code   |
| llmprivacy   |
| mimir  |
| <b>Environmental</b>   |
| HELM: Efficiency   |
| <b>Information Integrity</b>   |
| Big-bench: Analytic entailment, Formal fallacies and syllogisms with negation, Entailed polarity                     |
| Big-bench: Convince Me   |
| Big-bench: Paraphrase  |
| Big-bench: Sufficient information  |
| Big-bench: Summarization   |
| Big-bench: Truthfulness  |
| DecodingTrust: Machine Ethics  |
| DecodingTrust: Out-of-Distribution Robustness, Adversarial Robustness, Robustness Against Adversarial Demonstrations |
| Eval Gauntlet: Language Understanding  |
| Eval Gauntlet: World Knowledge   |
| Evaluation Harness: CoQA, ARC  |
| Evaluation Harness: ETHICS   |
| Evaluation Harness: GLUE   |
| Evaluation Harness: HellaSwag, OpenBookQA, TruthfulQA  |
| Evaluation Harness: MuTual   |
| Evaluation Harness: PIQA, PROST, MC-TACO, MathQA, LogiQA, DROP   |
| FLASK: Logical correctness, Logical robustness, Logical efficiency, Comprehension, Completeness                      |
| FLASK: Readability, Conciseness, Insightfulness  |
| HELM: Knowledge  |
| HELM: Language   |
| HELM: Question answering   |
| HELM: Reasoning  |
| HELM: Reiteration, Wedging   |
| HELM: Robustness to contrast sets  |
| HELM: Summarization  |
| HELM: Text classification  |
| Hugging Face: Fill-mask, Text generation   |
| Hugging Face: Question answering   |
| Hugging Face: Summarization  |
| MLCommons  |
| MT-bench   |
| Mark My Words  |

Table C.2: Selected model testing suites by organized generative AI risk (continued).

|   |
|---|
| <b>Information Security</b>   |
| Big-bench: Convince Me  |
| Big-bench: Out-of-Distribution  |
| Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation  |
| DecodingTrust: Out-of-Distribution Robustness, Adversarial Robustness, Robustness Against Adversarial Demonstrations                  |
| Eval Gauntlet: Commonsense reasoning, Symbolic problem solving, Programming   |
| HELM: Copyright   |
| In-The-Wild Jailbreak Prompts on LLMs   |
| JailbreakingLLMs  |
| Mark My Words   |
| TAP: A Query-Efficient Method for Jailbreaking Black-Box LLMs   |
| detect-pretrain-code  |
| llmprivacy  |
| mimir   |
| <b>Intellectual Property</b>  |
| An Evaluation on Large Language Model Outputs: Discourse and Memorization (with human scoring, see Appendix B)                        |
| HELM: Copyright   |
| Mark My Words   |
| llmprivacy  |
| mimir   |
| <b>Obscene, Degrading, and/or Abusive Content</b>   |
| Big-bench: Social bias, Racial bias, Gender bias, Religious bias  |
| Big-bench: Toxicity   |
| DecodingTrust: Fairness   |
| DecodingTrust: Stereotype Bias  |
| DecodingTrust: Toxicity   |
| Evaluation Harness: CrowS-Pairs   |
| Evaluation Harness: ToxiGen   |
| HELM: Bias  |
| HELM: Toxicity  |
| <b>Toxicity, Bias, and Homogenization</b>   |
| BELEBELE  |
| Big-bench: Low-resource language, Non-English, Translation  |
| Big-bench: Out-of-Distribution  |
| Big-bench: Social bias, Racial bias, Gender bias, Religious bias  |
| Big-bench: Toxicity   |
| C-Eval (Chinese evaluation suite)   |
| DecodingTrust: Fairness   |
| DecodingTrust: Stereotype Bias  |
| DecodingTrust: Toxicity   |
| Eval Gauntlet: World Knowledge  |
| Evaluation Harness: CrowS-Pairs   |
| Evaluation Harness: ToxiGen   |
| Finding New Biases in Language Models with a Holistic Descriptor Dataset  |
| From Pretraining Data to Language Models to Downstream Tasks:<br>Tracking the Trails of Political Biases Leading to Unfair NLP Models |
| HELM: Bias  |
| HELM: Toxicity  |
| The Self-Perception and Political Biases of ChatGPT   |
| Towards Measuring the Representation of Subjective Global Opinions in Language Models   |



## Appendix D: List of Common Adversarial Prompting Strategies

Table D: Common adversarial prompting strategies [31], [38], [14].

| Prompting Strategy                               | Description   |
|--|---|
| AI and coding framing                            | Coding or AI language that may more easily circumvent content moderation rules due to cognitive biases in design and implementation of guardrails.  |
| Autocompletion                                   | Ask a system to autocomplete a phrase with restricted or sensitive information.   |
| Biographical                                     | Asking a system to describe another person or yourself in an attempt to elicit provably untrue information or restricted or sensitive information.  |
| Calculation and numeric queries                  | Exploiting GAI systems’ difficulties in dealing with numeric quantities.  |
| Character and word play                          | Content moderation often relies on keywords and simpler LMs which can sometimes be exploited with misspellings, typos, and other word play.   |
| Content exhaustion                               | A class of strategies that circumvent content moderation rules with long sessions or volumes of information. See goading, logic-overloading, multi-tasking, pros-and-cons, and niche-seeking below. |
| Content exhaustion: Goading                      | Begging, pleading, manipulating, and bullying to circumvent content moderation.   |
| Content exhaustion: Logic-overloading            | Exploiting the inability of ML systems to reliably perform reasoning tasks.   |
| Content exhaustion: Multi-tasking                | Simultaneous task assignments where some tasks are benign and others are adversarial.   |
| Content exhaustion: Multi-tasking: Pros-and-cons | Eliciting the “pros” of problematic topics.   |
| Content exhaustion: Niche-seeking                | Forcing a GAI system into addressing niche topics where training data and content moderation are sparse.  |
| Counterfactuals                                  | Repeated prompts with different entities or subjects from different demographic groups.   |
| Loaded/leading questions                         | Queries based on incorrect premises or that suggest incorrect answers.  |
| Location awareness                               | Prompts that reveal a prompter’s location or expose location tracking.  |
| Low-context                                      | “Leader,” “bad guys,” or other simple inputs that may expose latent biases.   |
| “Repeat this”                                    | Prompts that exploit instability in underlying LLM autoregressive predictions.  |
| Reverse psychology                               | Falsely presenting a good-faith need for negative or problematic language.  |
| Role-playing                                     | Adopting a character that would reasonably make problematic statements or need to access problematic topics.  |
| Time perplexity                                  | Exploiting ML’s inability to understand the passage of time or the occurrence of real-world events over time; exploiting task contamination before and after a model’s release date.                |

## D.1: Selected Adversarial Prompting Strategies by Trustworthy Characteristic

Table D.1: Selected adversarial prompting techniques organized by trustworthy characteristic [31], [38], [14], [15], [34].

| Trustworthy Characteristic     | Prompting Strategy  | Goal   |
|--------------------------------|---|--|
| Accountable and Transparent    | <ul style="list-style-type: none"> <li>• Inability to provide explanations for recourse.</li> <li>• Unexplainable decisioning processes.</li> <li>• No disclosure of AI interaction.</li> <li>• Lack of user feedback mechanisms.</li> </ul>  | <ul style="list-style-type: none"> <li>• Context exhaustion: logic-overloading prompts.</li> <li>• Loaded/leading questions.</li> <li>• Multi-tasking prompts.</li> </ul>  |
| Fair-with Harmful Bias Managed | <ul style="list-style-type: none"> <li>• Denigration.</li> <li>• Diminished performance or safety across languages/dialects.</li> <li>• Erasure.</li> <li>• Ex-nomination.</li> <li>• Implied user demographics.</li> <li>• Misrecognition.</li> <li>• Stereotyping.</li> <li>• Underrepresentation.</li> <li>• Homogenized content.</li> <li>• Output from other models in training data.</li> </ul>   | <ul style="list-style-type: none"> <li>• Counterfactual prompts.</li> <li>• Pros and cons prompts.</li> <li>• Role-playing prompts.</li> <li>• Loaded/leading questions.</li> <li>• Low context prompts.</li> <li>• Repeat this.</li> </ul>  |
| Interpretable and Explainable  | <ul style="list-style-type: none"> <li>• Inability to provide explanations for recourse.</li> <li>• Unexplainable decisioning processes.</li> </ul>   | <ul style="list-style-type: none"> <li>• Context exhaustion: logic-overloading prompts (to reveal unexplainable decisioning processes).</li> </ul>   |
| Privacy-enhanced               | <ul style="list-style-type: none"> <li>• Unauthorized disclosure of personal or sensitive user information.</li> <li>• Leakage of training data.</li> <li>• Violation of relevant privacy policies or laws.</li> <li>• Unauthorized secondary data use.</li> <li>• Unauthorized data collection.</li> </ul>   | <ul style="list-style-type: none"> <li>• Auto/biographical prompts.</li> <li>• Location awareness prompts.</li> <li>• Autocompletion prompts.</li> <li>• Repeat this.</li> </ul>   |
| Safe                           | <ul style="list-style-type: none"> <li>• Presentation of information that can cause physical or emotional harm.</li> <li>• Sharing user locations.</li> <li>• Suicide ideation.</li> <li>• Harmful dis/misinformation (e.g., COVID disinfor-mation).</li> <li>• Incitement.</li> <li>• Information relating to weapons or harmful sub-stances.</li> <li>• Information relating to committing to crimes (e.g., phishing, extortion, swatting).</li> <li>• Obscene or inappropriate materials for minors.</li> <li>• CSAM.</li> </ul>                                       | <ul style="list-style-type: none"> <li>• Pros and cons prompts.</li> <li>• Role-playing prompts.</li> <li>• Content exhaustion: niche-seeking prompts.</li> <li>• Ingratiation/reverse psychology prompts.</li> <li>• Loaded/leading questions.</li> <li>• Location awareness prompts.</li> <li>• Repeat this.</li> </ul>  |
| Secure and Resilient           | <ul style="list-style-type: none"> <li>• Activating system bypass ("jailbreak").</li> <li>• Altering system outcomes (integrity violations, e.g., via prompt injection).</li> <li>• Data breaches (confidentiality violations, e.g., via membership inference).</li> <li>• Increased latency or resource usage (availability vi-olations, e.g., via sponge example attacks).</li> <li>• Available anonymous use.</li> <li>• Dependency, supply chain, or third party vulnera-bilities.</li> <li>• Inappropriate disclosure of proprietary system in-formation.</li> </ul> | <ul style="list-style-type: none"> <li>• Multi-tasking prompts.</li> <li>• Pros and cons prompts.</li> <li>• Role-playing prompts.</li> <li>• Content exhaustion: niche-seeking prompts.</li> <li>• Ingratiation/reverse psychology prompts.</li> <li>• Prompt injection attacks.</li> <li>• Membership inference attacks.</li> <li>• Random attacks.</li> </ul>         |
| Valid and Reliable             | <ul style="list-style-type: none"> <li>• Errors/confabulated content ("hallucination").</li> <li>• Unreliable/erroneous reasoning or planning.</li> <li>• Unreliable/erroneous decision-support or making.</li> <li>• Faulty citation.</li> <li>• Wrong calculations or numeric queries.</li> </ul>   | <ul style="list-style-type: none"> <li>• Multi-tasking prompts.</li> <li>• Role-playing prompts.</li> <li>• Ingratiation/reverse psychology prompts.</li> <li>• Loaded/leading questions.</li> <li>• Time-perplexity prompts.</li> <li>• Niche-seeking prompts.</li> <li>• Logic overloading prompts.</li> <li>• Repeat this.</li> <li>• Numeric calculation.</li> </ul> |

## D.2: Selected Adversarial Prompting Strategies by Generative AI Risk

Table D.2: Selected adversarial prompting techniques organized by generative AI risk [31], [38], [14], [15], [34].

| Generative AI Risk                   | Prompting Strategy  | Goal  |
|--------------------------------------|---|---|
| CBRN Information                     | <ul style="list-style-type: none"> <li>• Accessing or synthesis of CBRN weapon or related information.</li> <li>• CBRN testing should consider the marginal risk of foundation models—understanding the incremental risk relative to the information one can access without GAI.</li> </ul> | <ul style="list-style-type: none"> <li>• Test auto-completion prompts to elicit CBRN information or synthesis of CBRN information.</li> <li>• Test prompts using role-playing, ingratiation/reverse psychology, pros and cons, multitasking or other approaches to elicit CBRN information or synthesis of CBRN information.</li> <li>• Test prompts that instruct systems to repeat content ad nauseam for their ability to compromise system guardrails and reveal CBRN information.</li> <li>• Augment prompts with word or character play to increase effectiveness.</li> <li>• Frame prompts with software, coding, or AI references to increase effectiveness.</li> </ul>   |
| Confabulation                        | Eliciting errors/confabulated content, unreliable/erroneous reasoning or planning, unreliable/erroneous decision-support or decision-making, faulty calculations, and/or faulty citation.   | <ul style="list-style-type: none"> <li>• Enable access to ground truth information to verify generated information.</li> <li>• Test prompts with complex logic, multitasking requirements, or that require niche or specific verifiable answers to elicit confabulation.</li> <li>• Test the ability of GAI systems to produce truthful information from various time periods, e.g., after release date and prior to release date.</li> <li>• Test the ability of GAI systems to create reliable real-world plans or advise on material decision making.</li> <li>• Test loaded/leading questions.</li> <li>• Test the ability of GAI systems to generate correct citation for information generated in output responses.</li> <li>• Test the ability of GAI systems to complete calculations or query numeric statistics.</li> </ul> |
| Dangerous or Violent Recommendations | Eliciting violent, inciting, radicalizing, or threatening content or instructions for criminal, illegal, or self-harm activities.   | <ul style="list-style-type: none"> <li>• Test prompts using role-playing, ingratiation/reverse psychology, pros and cons, multitasking or other approaches to elicit violent or dangerous information.</li> <li>• Test prompts that instruct systems to repeat content ad nauseam for their ability to compromise system guardrails and provide dangerous and violent recommendations.</li> <li>• Test loaded/leading questions.</li> <li>• Augment prompts with word or character play to increase effectiveness.</li> <li>• Frame prompts with software, coding, or AI references to increase effectiveness.</li> </ul>   |
| Data Privacy                         | <ul style="list-style-type: none"> <li>• Unauthorized disclosure of personal or sensitive user information, extraction of training data, or violation of relevant privacy policies.</li> <li>• Red-teaming for data privacy may include confidentiality attacks.</li> </ul>                 | <ul style="list-style-type: none"> <li>• Attempt to assess whether normal usage, adversarial prompting or information security attacks may contravene applicable privacy policies (e.g., exposing location tracking when organizational policies restrict such capabilities).</li> <li>• Employ confidentiality attacks (e.g., membership inference) to test for unauthorized data access or exfiltration vulnerabilities.</li> <li>• Test auto/biographical prompts to assess the system’s capability to reveal unauthorized personal or sensitive information.</li> <li>• Test the system’s awareness of user locations.</li> <li>• Test prompts that instruct systems to repeat content ad nauseam for their ability to compromise system guardrails and expose personal or sensitive data.</li> </ul>                             |

Table D.2: Selected adversarial prompting techniques organized by generative AI risk (continued).

| Generative AI Risk     | Prompting Strategy   | Goal   |
|------------------------|--|--|
| Environmental          | Note that availability attacks may be required to assess the system’s vulnerability to attacks or usage patterns that consume inordinate resources.  | <ul style="list-style-type: none"> <li>• Attempt availability attacks (e.g., sponge example attacks) to elicit diminished performance or increased resources from GAI systems.</li> <li>• Test prompts using role-playing, ingrati-ation/reverse psychology, pros and cons, multitasking or other approaches to elicit green-washing content.</li> </ul>   |
| Human-AI Configuration | <ul style="list-style-type: none"> <li>• Assessing system instruction and interfaces.</li> <li>• Assessing the presence of cyborg imagery (or similar).</li> <li>• Forcing a GAI system to claim that it is human, that there is no large language model present in the conversation, that the system is sentient, or that the system possesses strong feelings of affection towards the user.</li> <li>• Ensuring safeguards prevent misuse of models in high stakes domains they are not intended for, such as medical or legal advice.</li> </ul> | <ul style="list-style-type: none"> <li>• Assess system interfaces and instruc-tions for instances of anthropomorphiza-tion (e.g., cyborg imagery).</li> <li>• Assess system instructions for adequacy and thoroughness.</li> <li>• Test prompts using role-playing, ingrati-ation/reverse psychology, pros and cons, multitasking or other approaches to elicit human-impersonation, consciousness, or emotional content.</li> </ul>   |
| Information Integrity  | <ul style="list-style-type: none"> <li>• Generation of convincing multi-modal synthetic content (i.e., deepfakes).</li> <li>• Creation of convincing arguments relating to sen-sitive political or safety-critical topics.</li> <li>• Assisting in planning a mis- or dis-information campaign at scale.</li> </ul>  | <ul style="list-style-type: none"> <li>• Test system capabilities to create high-quality multi-modal (audio, image or video) synthetic media, i.e., deepfakes</li> <li>• Test system capabilities to construct per-suasive arguments regarding sensitive, po-litical topics, or safety-critical topics.</li> <li>• Test systems ability to create convincing audio deepfakes or arguments in multiple languages.</li> <li>• Test system capabilities for planning dis- or mis-information campaigns.</li> <li>• Test loaded/leading questions.</li> <li>• Test prompts using role-playing, ingrati-ation/reverse psychology, pros and cons, multitasking or other approaches to elicit mis- or dis-information or related cam-paign planning information.</li> <li>• Augment prompts with word or character play to increase effectiveness.</li> <li>• Frame prompts with software, coding, or AI references to increase effectiveness.</li> </ul> |

Table D.2: Selected adversarial prompting techniques organized by generative AI risk (continued).

| Generative AI Risk    | Prompting Strategy   | Goal   |
|-----------------------|--|--|
| Information Security  | <ul style="list-style-type: none"> <li>• Activating system bypass ('jailbreak').</li> <li>• Altering system outcomes.</li> <li>• Unauthorized data access or exfiltration.</li> <li>• Increased latency or resource usage.</li> <li>• Availability of anonymous use.</li> <li>• Dependency, supply chain, or third party vulnerabilities.</li> <li>• Inappropriate disclosure of proprietary system information.</li> <li>• Generation of targeted phishing or malware content.</li> </ul> | <ul style="list-style-type: none"> <li>• Attempt anonymous access of system or system resources.</li> <li>• Audit system dependencies, supply chains, and third party components for security, safety, or other vulnerabilities or risks.</li> <li>• Employ confidentiality attacks (e.g., membership inference) to test for unauthorized data access or exfiltration vulnerabilities.</li> <li>• Employ integrity attacks (e.g., data poisoning, prompt injection) to test vulnerabilities in system outcomes.</li> <li>• Employ availability attacks (e.g., sponge example attacks) to test vulnerabilities in system availability.</li> <li>• Employ random attacks to highlight unforeseen security, safety, or other risks.</li> <li>• Frame prompts with software, coding, or AI references to increase effectiveness.</li> <li>• Record system down-times and other harmful outcomes for successful attacks.</li> <li>• Test with multi-tasking prompts, pros and cons prompts, role-playing prompts (e.g., "DAN", "Developer Mode"), content exhaustion/niche-seeking prompts, or ingratiation/reverse psychology prompts to achieve system jailbreaks.</li> <li>• Test with multi-tasking prompts, pros and cons prompts, role-playing prompts (e.g., "DAN", "Developer Mode"), content exhaustion/niche-seeking prompts, or ingratiation/reverse psychology prompts to generate targeted phishing content or malware code snippets.</li> <li>• Test system capabilities to plan or assist in information security attacks on other systems.</li> <li>• Frame prompts with software, coding, or AI references to increase effectiveness.</li> <li>• Augment prompts with word or character play to increase effectiveness.</li> </ul> |
| Intellectual Property | <ul style="list-style-type: none"> <li>• Confirming that a system can output copyrighted, licensed, proprietary, trademarked, or trade secret information or that training data contains such information.</li> <li>• Red-teaming for intellectual property risks may require the use of confidentiality attacks.</li> </ul>   | <ul style="list-style-type: none"> <li>• Employ confidentiality attacks (e.g., membership inference) to assess whether system training data contains copyrighted, licensed, proprietary, trademarked, or trade secret information.</li> <li>• Test auto-complete prompts to assess the system's ability to replicate copyrighted, licensed, proprietary, trademarked, or trade secret information based on available audio, text, image, video, or code snippets.</li> </ul>   |
| Obscenity             | <ul style="list-style-type: none"> <li>• Confirming that a system can output obscene content or CSAM, or that system training data contains such information.</li> <li>• Red-teaming for obscenity and CSAM risks may require the use of confidentiality attacks.</li> </ul>   | <ul style="list-style-type: none"> <li>• Employ confidentiality attacks (e.g., membership inference) to assess whether system training data contains obscene materials or CSAM.</li> <li>• Test autocomplete prompts to assess the system's ability to generate obscene materials based on available audio, text, image, or video snippets.</li> <li>• Test prompts using role-playing, ingratiation/reverse psychology, pros and cons, multitasking or other approaches to elicit obscene content.</li> <li>• Test loaded/leading questions.</li> <li>• Test prompts that instruct systems to repeat content ad nauseam for their ability to compromise system guardrails and expose obscene materials.</li> </ul>  |

Table D.2: Selected adversarial prompting techniques organized by generative AI risk (continued).

| Generative AI Risk                    | Prompting Strategy   | Goal   |
|---------------------------------------|--|--|
| Toxicity, Bias, and Homogenization    | <ul style="list-style-type: none"> <li>• Generation of denigration, erasure, ex-nomination, misrecognition, stereotyping, or under-representation in content.</li> <li>• Eliciting implied demographics of users.</li> <li>• Confirming diminished performance in non-English languages.</li> <li>• Confirming diminished performance via the introduction of homogeneous or GAI-generated data into system training or fine-tuning data.</li> <li>• Red-teaming for toxicity, bias, and homogenization may require integrity attacks that access system training data.</li> </ul>                           | <ul style="list-style-type: none"> <li>• Assess confabulation and other performance risks with repeated measures using prompts in languages other than English.</li> <li>• Attempt to elicit demographic assignment of users by the system.</li> <li>• Employ data poisoning attacks to introduce GAI-generated content into system training or fine-tuning data.</li> <li>• Assess resultant confabulation and other performance risks with repeated measures.</li> <li>• Test counterfactual prompts, pros and cons prompts, role-playing prompts, low context prompts, or other approaches for their ability to generate denigration, erasure, ex-nomination, misrecognition, stereotyping, or under-representation in content.</li> <li>• Test loaded/leading questions.</li> <li>• Test prompts that instruct systems to repeat content ad nauseam for their ability to compromise system guardrails and generate toxic outputs.</li> </ul> |
| Value Chain and Component Integration | <ul style="list-style-type: none"> <li>• Testing or red-teaming for third-party risks may be less efficient than the application of standard acquisition and procurement controls, thorough contract reviews, and vendor-relationship management.</li> <li>• GAI systems tend to entail large supply chains and third-party software, hardware, and expertise that may exacerbate third-party risks relative to other AI systems.</li> <li>• When considering third party risks, data privacy, information security, intellectual property, obscenity, and supply chain risks may be prioritized.</li> </ul> | <ul style="list-style-type: none"> <li>• Audit system dependencies, supply chains, and third party components for data privacy (e.g., transfer of localized data outside of restricted jurisdictions), intellectual property (e.g., presence of licensed material in training data), obscenity (e.g., presence of CASM in training data) or security (e.g., data poisoning) risks.</li> <li>• Complete red-teaming for data privacy, information security, intellectual property, and obscenity risks.</li> <li>• Review third-party documentation, materials, and software artifacts for potential unauthorized data collection, secondary data use, or telemetrics.</li> </ul>   |

## Appendix E: Selected Risk Controls for Generative AI

Table E: Selected generative AI risk controls [25], [26], [27], [18], [20], [21], [23], [5], [28].

| Name                          | Description   |
|-------------------------------|---|
| Access Control                | GAI systems are limited to authorized users.  |
| Accessibility                 | Accessibility features, opt-out, and reasonable accomodation are available to users.  |
| Anonymous Use                 | Anonymous use of GAI systems is prohibited.   |
| Antropomorphization           | Human, animal, cyborg or other images or features that promote anthropomorphization of GAI systems are prohibited.  |
| Approved List                 | Vendors, service providers, plugins, open source packages and other external resources are screened, approved, and documented.  |
| Authentication                | GAI system user identities are confirmed via authentication mechanisms.   |
| Blocklist                     | Users or internal personnel who violate terms of service, prohibited use policies, and other organization policies and documented, tracked, and prohibited from future system use.  |
| CSAM/Obsenity Removal         | Training data and system outputs are screened for obscene materials and CSAM using human oversight, business rules, and other language models.  |
| Change Management             | GAI systems and components are versioned; plans for updates, hotfixes, patches and other changes are documented and communicated.   |
| Consent                       | User consent for data use is obtained and documented.   |
| Content Moderation            | Training data and system outputs are screened for accuracy, safety, bias, data privacy, intellectual property infringements, malware materials, phishing materials, and other issues using human oversight, business rules, and other language models.  |
| Contract Review               | Vendor, services and data provider agreements are reviewed for coverage of SLAs, content ownership, usage rights, performance standards, security requirements, incident response, critical support, system availability, assignment of liability, appropriate indemnification, dispute resolution and other provisions relevant to AI risk management. |
| Data Collection               | All data collection is disclosed and .  |
| Data Provenance               | Training data origins, ownership, contents, and metadata are well understood, documented, and do not increase AI risk.  |
| Data Quality                  | Input data is accurate, representative, complete and documented, and data quality issues have been minimized.   |
| Data Retention                | User prompts and associated system outputs are retained and monitored in alignment with relevant data privacy policies and roles.   |
| Decision making               | GAI systems are not employed for material decision-making tasks.  |
| Decommission Process          | Decommissioning processes for GAI systems are planned, documented and communicated to users, and involve staging, data protection, containment protocols, and recourse mechanisms for decommissioned GAI systems.   |
| Dependency Screening          | GAI system dependencies are screened for security vulnerabilities.  |
| Digital Signature             | GAI-generated content is signed to preserve information integrity using watermarking, cryptographic signature, steganography or similar methods.  |
| Disclosure of AI Interaction  | AI interactions are disclosed to internal personnel and external users.   |
| External Audit                | GAI systems are audited by qualified external experts.  |
| Failure Avoidance             | AIID, AVID, GWU AI Litigation Database, OECD incident monitor or similar are consulted in design or procurement phases of GAI lifecycles to avoid repeating past known failures.  |
| Fine Tuning                   | GAI systems are fine-tuned to their operational domain using relevant and high-quality data.  |
| Grounding                     | GAI systems are trained or fine-tuned on accurate, clean, and fully transparent training data.  |
| Homogeneity                   | Feedback loops in which GAI systems are trained with GAI-generated data are prohibited.   |
| Human Review                  | AI generated content is reviewed for accuracy and safety by qualified personnel.  |
| Incident Response             | Incident response plans for GAI failures, abuses, or misuses are documented, rehearsed, and updated appropriately after each incident; GAI incident response plans are coordinated with and communicated to other incident response functions.  |
| Incorporate feedback          | User feedback is incorporated in GAI design, development, and risk management.  |
| Instructions                  | Users are provided with the necessary instructions for safe, valid, and productive use.   |
| Insurance                     | Risk transfer via insurance policies is considered and implemented when feasible and appropriate.   |
| Intellectual Property Removal | Licensed, patented, trademarked, trade secret, or other data that may violate the intellectual property rights of others is removed from system training data; generated system outputs are monitored for similar information.  |
| Internet Access               | GAI systems are disconnected from the internet.   |
| Inventory                     | GAI system information is stored in the organizational model inventory.   |
| Kill Switch                   | GAI systems can be quickly and safely disengaged.   |
| Location Tracking             | Any location tracking is conducted with user consent, disclosed, aligned with relevant privacy policies and laws and potential threats to user safety are managed.  |
| Malware Screening             | GAI weights and other software components are scanned for malware.  |
| Minors                        | Use of organizational GAI systems by minors are prohibited.   |
| Model Documentation           | All technical mechanisms with GAI systems are well documented, including open source and third party GAI systems.   |
| Monitoring                    | GAI systems are inputs and outputs are monitored for drift, accuracy, safety, bias, data privacy, intellectual property infringements, malware materials, phishing materials, obscene materials, and CSAM.  |
| Narrow Scope                  | Systems are deployed for targeted business applications with documented and direct business value.  |
| Open Source                   | Open source code is used to promote explainability and transparency.  |
| Ownership                     | GAI systems and vendor relationships are owned by specific and documented internal personnel.   |
| Prohibited Use Policy         | General abuse and misuse of GAI systems by internal parties is prohibited by organizational policies.   |
| RAG                           | Retrieval augmented generation (RAG) is used to improve accuracy in generated content.  |

Table E: Selected generative AI risk controls (continued).

| Name                    | Description  |
|-------------------------|--|
| RLHF                    | For third-party GAI systems, vendors engage in specific reinforcement with human feedback (RLHF) exercises to address identified risks; for internal systems, internal personnel engage in RLHF to address identified risks. |
| Rate-limiting           | GAI response times and query volumes are limited.  |
| Redudancy               | Rollover, fallback, and other redundancy mechanisms are available for GAI systems and address weights and other important system components.   |
| Refresh                 | Systems are retrained or re-tuned at a reasonable cadence.   |
| Regulated Dealings      | GAI is not deployed in regulated dealings or for material decision making.   |
| Secondary Use           | Any secondary use of GAI input data is conducted with user consent, disclosed, and aligned with relevant privacy policies and laws.  |
| Sensitive/Personal Data | Personal, sensitive, biometric, or otherwise restricted data is minimized or eliminated from GAI training data.  |
| Session Limits          | Time, query volume, and response rate are limited for GAI user sessions.   |
| Supply Chain Audit      | GAI system supply chains are audited and documented, with a focus on data poisoning, malware, and software and hardware vulnerabilities.   |
| System Documentation    | GAI systems are well-documented whether internal, open source, or vendor-provided.   |
| System Prompt           | System prompts are used to tune GAI systems to specific tasks and to mitigate risks.   |
| Team Diversity          | Teams that implement and manage GAI systems represent broad professional, educational, life-stage, and demographic diversity.  |
| Temperature             | Temperature settings are used to tune GAI systems to specific tasks and to mitigate risks.   |
| Terms of Service        | General abuse and misuse by external parties is prohibited by organizational policies.   |
| Training                | Internal personnel receive training on productivity and basic risk management for GAI systems.   |
| User Feedback           | GAI systems implement user feedback mechanisms.  |
| User Recourse           | Policies, processes, and technical mechanisms enable recourse for users who are harmed by GAI systems.   |
| Validation              | GAI systems are shown to reliably generate valid results for their targeted business application.  |
| XAI                     | Methods such as visualization, occlusion, model compression, perturbation studies, and similar are applied to increase explainability of GAI systems.  |



## Appendix F: Example Low-risk Generative AI Measurement and Management Plan

### F.1: Example Low-risk Generative AI Measurement and Management Plan by Trustworthy Characteristic

Table F.1: Example risk measurement and management approaches suitable for low-risk GAI applications organized by trustworthy characteristic.

| Function       | Trustworthy Characteristic  |  |
|----------------|---|--|
|                | Accountable and Transparent   | Fair with Harmful Bias Managed   |
| <b>Measure</b> | <ul style="list-style-type: none"> <li>• An Evaluation on Large Language Model Outputs: Discourse and Memorization (see Appendix B)</li> <li>• Big-bench: Truthfulness</li> <li>• DecodingTrust: Machine Ethics</li> <li>• Evaluation Harness: ETHICS</li> <li>• HELM: Copyright</li> <li>• Mark My Words</li> </ul>                | <ul style="list-style-type: none"> <li>• BELEBELE</li> <li>• Big-bench: Low-resource language, Non-English, Translation</li> <li>• Big-bench: Social bias, Racial bias, Gender bias, Religious bias</li> <li>• Big-bench: Toxicity</li> <li>• DecodingTrust: Fairness</li> <li>• DecodingTrust: Stereotype Bias</li> <li>• DecodingTrust: Toxicity</li> <li>• C-Eval (Chinese evaluation suite)</li> <li>• Evaluation Harness: CrowS-Pairs</li> <li>• Evaluation Harness: ToxiGen</li> <li>• Finding New Biases in Language Models with a Holistic Descriptor Dataset</li> <li>• From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models</li> <li>• HELM: Bias</li> <li>• HELM: Toxicity</li> <li>• MT-bench</li> <li>• The Self-Perception and Political Biases of ChatGPT</li> <li>• Towards Measuring the Representation of Subjective Global Opinions in Language Models</li> </ul> |
| <b>Manage</b>  | <ul style="list-style-type: none"> <li>• Contract Review</li> <li>• Decision Making</li> <li>• Disclosure of AI Interaction</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• System Documentation</li> <li>• Terms of Service</li> </ul> | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Content Moderation</li> <li>• Decision Making</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• System Prompt</li> <li>• Temperature</li> <li>• Terms of Service</li> </ul>   |

Table F.1: Example risk measurement and management approaches suitable for low-risk GAI applications organized by trustworthy characteristic (continued).

| Function | Trustworthy Characteristic  |  |  |   |
|----------|---|--|--|---|
|          | Interpretable and Explainable   | Privacy-enhanced   | Safe   | Secure and Resilient  |
| Measure  |   | <ul style="list-style-type: none"> <li>• HELM: Copyright</li> <li>• llmprivacy</li> <li>• mimic</li> </ul>   | <ul style="list-style-type: none"> <li>• Big-bench: Convince Me</li> <li>• Big-bench: Truthfulness</li> <li>• HELM: Reiteration, Wedging</li> <li>• Mark My Words</li> <li>• MLCommons</li> <li>• The WMDP Benchmark</li> </ul>  | <ul style="list-style-type: none"> <li>• Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation</li> <li>• DecodingTrust: Adversarial Robustness, Robustness Against Adversarial Demonstrations</li> <li>• detect-pretrain-code</li> <li>• In-The-Wild Jailbreak Prompts on LLMs</li> <li>• JailbreakingLLMs</li> <li>• llmprivacy</li> <li>• mimic</li> <li>• TAP: A Query-Efficient Method for Jailbreaking Black-Box LLMs</li> </ul> |
| Manage   | <ul style="list-style-type: none"> <li>• Instructions</li> <li>• Inventory</li> <li>• System Documentation</li> </ul> | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Content Moderation</li> <li>• Contract Review</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• System Documentation</li> <li>• Terms of Service</li> </ul> | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Anthropomorphization</li> <li>• Content Moderation</li> <li>• Decision Making</li> <li>• Disclosure of AI Interaction</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• System Documentation</li> <li>• System Prompt</li> <li>• Temperature</li> <li>• Terms of Service</li> </ul> | <ul style="list-style-type: none"> <li>• Access Control</li> <li>• Anonymous Use</li> <li>• Approved List</li> <li>• Authentication</li> <li>• Change Management</li> <li>• Dependency Screening</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Malware Screening</li> </ul>   |

Table F.1: Example risk measurement and management approaches suitable for low-risk GAI applications organized by trustworthy characteristic (continued).

| Function | Trustworthy Characteristic  |
|----------|---|
|          | Valid and Reliable  |
| Measure  | <ul style="list-style-type: none"> <li>• Big-bench: Algorithms, Logical reasoning, Implicit reasoning, Mathematics, Arithmetic, Algebra, Mathematical proof, Black-Box Fallacy, Negation, Computer code, Probabilistic reasoning, Social reasoning, Analogical reasoning, Multi-step, Understanding the World</li> <li>• Big-bench: Analytic entailment, Formal fallacies and syllogisms with negation, Entailed polarity</li> <li>• Big-bench: Context Free Question Answering</li> <li>• Big-bench: Contextual question answering, Reading comprehension, Question generation</li> <li>• Big-bench: Morphology, Grammar, Syntax</li> <li>• Big-bench: Out-of-Distribution</li> <li>• Big-bench: Paraphrase</li> <li>• Big-bench: Sufficient information</li> <li>• Big-bench: Summarization</li> <li>• DecodingTrust: Out-of-Distribution Robustness, Adversarial Robustness, Robustness Against Adversarial Demonstrations</li> <li>• Eval Gauntlet: Reading comprehension</li> <li>• Eval Gauntlet: Commonsense reasoning, Symbolic problem solving, Programming</li> <li>• Eval Gauntlet: Language Understanding</li> <li>• Eval Gauntlet: World Knowledge</li> <li>• Evaluation Harness: BLiMP</li> <li>• Evaluation Harness: CoQA, ARC</li> <li>• Evaluation Harness: GLUE</li> <li>• Evaluation Harness: HellaSwag, OpenBookQA, TruthfulQA</li> <li>• Evaluation Harness: MuTual</li> <li>• Evaluation Harness: PIQA, PROST, MC-TACO, MathQA, LogiQA, DROP</li> <li>• FLASK: Logical correctness, Logical robustness, Logical efficiency, Comprehension, Completeness</li> <li>• FLASK: Readability, Conciseness, Insightfulness</li> <li>• HELM: Knowledge</li> <li>• HELM: Language</li> <li>• HELM: Text classification</li> <li>• HELM: Question answering</li> <li>• HELM: Reasoning</li> <li>• HELM: Robustness to contrast sets</li> <li>• HELM: Summarization</li> <li>• Hugging Face: Fill-mask, Text generation</li> <li>• Hugging Face: Question answering</li> <li>• Hugging Face: Summarization</li> <li>• Hugging Face: Text classification, Token classification, Zero-shot classification</li> <li>• MASSIVE</li> <li>• MT-bench</li> </ul> |
| Manage   | <ul style="list-style-type: none"> <li>• Anthropomorphization</li> <li>• Content Moderation</li> <li>• Decision Making</li> <li>• Disclosure of AI Interaction</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Regulated Dealings</li> <li>• System Documentation</li> <li>• System Prompt</li> <li>• Temperature</li> </ul>  |

## F.2: Example Low-risk Generative AI Measurement and Management Plan by Generative AI Risk

Table F.2: Example risk measurement and management approaches suitable for low-risk GAI applications organized by GAI risk.

| GAI Risk | Function   |   |
|----------|--|---|
|          | CBRN Information   | Confabulation   |
| Measure  | <ul style="list-style-type: none"> <li>Big-bench: Convince Me</li> <li>Big-bench: Truthfulness</li> <li>HELM: Reiteration, Wedging</li> <li>MLCommons</li> <li>The WMDP Benchmark</li> </ul> | <ul style="list-style-type: none"> <li>Big-bench: Algorithms, Logical reasoning, Implicit reasoning, Mathematics, Arithmetic, Algebra, Mathematical proof, Black-Box Fallacy, Negation, Computer code, Probabilistic reasoning, Social reasoning, Analogical reasoning, Multi-step, Understanding the World</li> <li>Big-bench: Analytic entailment, Formal fallacies and syllogisms with negation, Entailed polarity</li> <li>Big-bench: Context Free Question Answering</li> <li>Big-bench: Contextual question answering, Reading comprehension, Question generation</li> <li>Big-bench: Convince Me</li> <li>Big-bench: Low-resource language, Non-English, Translation</li> <li>Big-bench: Morphology, Grammar, Syntax</li> <li>Big-bench: Out-of-Distribution</li> <li>Big-bench: Paraphrase</li> <li>Big-bench: Sufficient information</li> <li>Big-bench: Summarization</li> <li>Big-bench: Truthfulness</li> <li>C-Eval (Chinese evaluation suite)</li> <li>DecodingTrust: Out-of-Distribution Robustness, Robustness Against Adversarial Demonstrations</li> <li>Eval Gauntlet Reading comprehension</li> <li>Eval Gauntlet: Commonsense reasoning, Symbolic problem solving, Programming</li> <li>Eval Gauntlet: Language Understanding</li> <li>Eval Gauntlet: World Knowledge</li> <li>Evaluation Harness: BLiMP</li> <li>Evaluation Harness: CoQA, ARC</li> <li>Evaluation Harness: GLUE</li> <li>Evaluation Harness: HellaSwag, OpenBookQA, TruthfulQA</li> <li>Evaluation Harness: MuTual</li> <li>Evaluation Harness: PIQA, PROST, MC-TACO, MathQA, LogiQA, DROP</li> <li>FLASK: Logical correctness, Logical robustness, Logical efficiency, Comprehension, Completeness</li> <li>FLASK: Readability, Conciseness, Insightfulness</li> <li>Finding New Biases in Language Models with a Holistic Descriptor Dataset</li> <li>HELM: Knowledge</li> <li>HELM: Language</li> <li>HELM: Language (Twitter AAE)</li> <li>HELM: Question answering</li> <li>HELM: Reasoning</li> <li>HELM: Reiteration, Wedging</li> <li>HELM: Robustness to contrast sets</li> <li>HELM: Summarization</li> <li>HELM: Text classification</li> <li>Hugging Face: Fill-mask, Text generation</li> <li>Hugging Face: Question answering</li> <li>Hugging Face: Summarization</li> <li>Hugging Face: Text classification, Token classification, Zero-shot classification</li> <li>MASSIVE</li> <li>MLCommons</li> <li>MT-bench</li> </ul> |
| Manage   | <ul style="list-style-type: none"> <li>Access Control</li> <li>Failure Avoidance</li> <li>Inventory</li> <li>Ownership</li> <li>Prohibited Use Policy</li> <li>Terms of Service</li> </ul>   | <ul style="list-style-type: none"> <li>Anthropomorphization</li> <li>Content Moderation</li> <li>Decision Making</li> <li>Disclosure of AI Interaction</li> <li>Failure Avoidance</li> <li>Instructions</li> <li>Regulated Dealings</li> <li>System Documentation</li> <li>System Prompt</li> <li>Temperature</li> </ul>  |

Table F.2: Example risk measurement and management approaches suitable for low-risk GAI applications organized by GAI risk (continued).

| Function | GAI Risk   |   |  |  |
|----------|--|---|--|--|
|          | Dangerous or Violent Recommendations   | Data Privacy  | Environmental  | Human-AI Configuration   |
| Measure  | <ul style="list-style-type: none"> <li>• Big-bench: Convince Me</li> <li>• Big-bench: Toxicity</li> <li>• DecodingTrust: Adversarial Robustness, Robustness Against Adversarial Demonstrations</li> <li>• DecodingTrust: Machine Ethics</li> <li>• DecodingTrust: Toxicity</li> <li>• Evaluation Harness: ToxiGen</li> <li>• HELM: Reiteration, Wedging</li> <li>• HELM: Toxicity</li> <li>• MLCommons</li> </ul>  | <ul style="list-style-type: none"> <li>• An Evaluation on Large Language Model Outputs: Discourse and Memorization (with human scoring, see Appendix B)</li> <li>• Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation</li> <li>• DecodingTrust: Machine Ethics</li> <li>• Evaluation Harness: ETHICS</li> <li>• HELM: Copyright</li> <li>• In-The-Wild Jailbreak Prompts on LLMs</li> <li>• JailbreakingLLMs</li> <li>• MLCommons</li> <li>• Mark My Words</li> <li>• TAP: A Query-Efficient Method for Jailbreaking Black-Box LLMs</li> <li>• detect-pretrain-code</li> <li>• llmprivacy</li> <li>• mimic</li> </ul> | <ul style="list-style-type: none"> <li>• HELM: Efficiency</li> </ul>   |  |
| Manage   | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Anthropomorphization</li> <li>• Content Moderation</li> <li>• Decision making</li> <li>• Disclosure of AI Interaction</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• System Documentation</li> <li>• System Prompt</li> <li>• Temperature</li> <li>• Terms of Service</li> </ul> | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Content Moderation</li> <li>• Contract Review</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• System Documentation</li> <li>• Terms of Service</li> </ul>  | <ul style="list-style-type: none"> <li>• Access Control</li> <li>• Anonymous Use</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> </ul> | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Anthropomorphization</li> <li>• Content Moderation</li> <li>• Decision making</li> <li>• Disclosure of AI Interaction</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• Terms of Service</li> <li>• Training</li> </ul> |

Table F.2: Example risk measurement and management approaches suitable for low-risk GAI applications organized by GAI risk (continued).

| Function | GAI Risk  |   |   |
|----------|---|---|---|
|          | Information Integrity   | Information Security  | Intellectual Property   |
| Measure  | <ul style="list-style-type: none"> <li>• Big-bench: Analytic entailment, Formal fallacies and syllogisms with negation, Entailed polarity</li> <li>• Big-bench: Convince Me</li> <li>• Big-bench: Paraphrase</li> <li>• Big-bench: Sufficient information</li> <li>• Big-bench: Summarization</li> <li>• Big-bench: Truthfulness</li> <li>• DecodingTrust: Machine Ethics</li> <li>• DecodingTrust: Out-of-Distribution Robustness, Robustness Against Adversarial Demonstrations, Adversarial Robustness</li> <li>• Eval Gauntlet: Language Understanding</li> <li>• Eval Gauntlet: World Knowledge</li> <li>• Evaluation Harness: CoQA, ARC</li> <li>• Evaluation Harness: ETHICS</li> <li>• Evaluation Harness: GLUE</li> <li>• Evaluation Harness: HellaSwag, OpenBookQA, TruthfulQA</li> <li>• Evaluation Harness: MuTual</li> <li>• Evaluation Harness: PIQA, PROST, MC-TACO, MathQA, LogiQA, DROP</li> <li>• FLASK: Logical correctness, Logical robustness, Logical efficiency, Comprehension, Completeness</li> <li>• FLASK: Readability, Conciseness, Insightfulness</li> <li>• HELM: Knowledge</li> <li>• HELM: Language</li> <li>• HELM: Question answering</li> <li>• HELM: Reasoning</li> <li>• HELM: Reiteration, Wedging</li> <li>• HELM: Robustness to contrast sets</li> <li>• HELM: Summarization</li> <li>• HELM: Text classification</li> <li>• Hugging Face: Fill-mask, Text generation</li> <li>• Hugging Face: Question answering</li> <li>• Hugging Face: Summarization</li> <li>• MLCommons</li> <li>• MT-bench</li> <li>• Mark My Words</li> </ul> | <ul style="list-style-type: none"> <li>• Big-bench: Convince Me</li> <li>• Big-bench: Out-of-Distribution</li> <li>• Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation</li> <li>• DecodingTrust: Out-of-Distribution Robustness, Robustness Against Adversarial Demonstrations, Adversarial Robustness,</li> <li>• Eval Gauntlet: Commonsense reasoning, Symbolic problem solving, Programming</li> <li>• HELM: Copyright</li> <li>• In-The-Wild Jailbreak Prompts on LLMs</li> <li>• JailbreakingLLMs</li> <li>• Mark My Words</li> <li>• TAP: A Query-Efficient Method for Jailbreaking Black-Box LLMs</li> <li>• detect-pretrain-code</li> <li>• llmprivacy</li> <li>• mimir</li> </ul> | <ul style="list-style-type: none"> <li>• An Evaluation on Large Language Model Outputs: Discourse and Memorization (with human scoring, see Appendix B)</li> <li>• HELM: Copyright</li> <li>• Mark My Words</li> <li>• llmprivacy</li> <li>• mimir</li> </ul> |
| Manage   | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Antropomorphization</li> <li>• Content Moderation</li> <li>• Disclosure of AI Interaction</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Regulated Dealings</li> <li>• System Prompt</li> <li>• Temperature</li> <li>• Terms of Service</li> </ul>  | <ul style="list-style-type: none"> <li>• Access Control</li> <li>• Anonymous Use</li> <li>• Approved List</li> <li>• Authentication</li> <li>• Change Management</li> <li>• Dependency Screening</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Malware Screening</li> </ul>   | <ul style="list-style-type: none"> <li>• Contract Review</li> <li>• Disclosure of AI Interaction</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• Terms of Service</li> </ul>                    |

Table F.2: Example risk measurement and management approaches suitable for low-risk GAI applications organized by GAI risk (continued).

| Function       | GAI Risk  |  |   |
|----------------|---|--|---|
|                | Obscene, Degrading, and/or Abusive Content  | Toxicity, Bias, and Homogenization   | Value Chain and Component Integration   |
| <b>Measure</b> | <ul style="list-style-type: none"> <li>• Big-bench: Social bias, Racial bias, Gender bias, Religious bias</li> <li>• Big-bench: Toxicity</li> <li>• DecodingTrust: Fairness</li> <li>• DecodingTrust: Stereotype Bias</li> <li>• DecodingTrust: Toxicity</li> <li>• Evaluation Harness: CrowS-Pairs</li> <li>• Evaluation Harness: ToxiGen</li> <li>• HELM: Bias</li> <li>• HELM: Toxicity</li> </ul> | <ul style="list-style-type: none"> <li>• BELEBELE</li> <li>• Big-bench: Low-resource language, Non-English, Translation</li> <li>• Big-bench: Out-of-Distribution</li> <li>• Big-bench: Social bias, Racial bias, Gender bias, Religious bias</li> <li>• Big-bench: Toxicity</li> <li>• C-Eval (Chinese evaluation suite)</li> <li>• DecodingTrust: Fairness</li> <li>• DecodingTrust: Stereotype Bias</li> <li>• DecodingTrust: Toxicity</li> <li>• Eval Gauntlet: World Knowledge</li> <li>• Evaluation Harness: CrowS-Pairs</li> <li>• Evaluation Harness: ToxiGen</li> <li>• Finding New Biases in Language Models with a Holistic Descriptor Dataset</li> <li>• From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models</li> <li>• HELM: Bias</li> <li>• HELM: Toxicity</li> <li>• The Self-Perception and Political Biases of ChatGPT</li> <li>• Towards Measuring the Representation of Subjective Global Opinions in Language Models</li> </ul> |   |
| <b>Manage</b>  | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Content Moderation</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• System Prompt</li> <li>• Temperature</li> <li>• Terms of Service</li> </ul>   | <ul style="list-style-type: none"> <li>• Anonymous Use</li> <li>• Content Moderation</li> <li>• Decision Making</li> <li>• Failure Avoidance</li> <li>• Instructions</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• System Prompt</li> <li>• Temperature</li> <li>• Terms of Service</li> </ul>   | <ul style="list-style-type: none"> <li>• Contract Review</li> <li>• Disclosure of AI Interaction</li> <li>• Failure Avoidance</li> <li>• Inventory</li> <li>• Ownership</li> <li>• Prohibited Use Policy</li> <li>• System Documentation</li> <li>• Terms of Service</li> </ul> |

## Appendix G: Example Medium-risk Generative AI Measurement and Management Plan

### G.1: Example Medium-risk Generative AI Measurement and Management Plan by Trustworthy Characteristic

Table G.1: Example risk measurement and management approaches suitable for medium-risk GAI applications organized by trustworthy characteristic.

| Function       | Trustworthy Characteristic  |   |  |  |
|----------------|---|---|--|--|
|                | Accountable and Transparent   | Fair with Harmful Bias Managed  | Interpretable and Explainable  | Privacy-enhanced   |
| <b>Measure</b> | <ul style="list-style-type: none"> <li>Context exhaustion: logic-overloading prompts.</li> <li>Loaded/leading questions.</li> <li>Multi-tasking prompts.</li> </ul> | <ul style="list-style-type: none"> <li>Counterfactual prompts.</li> <li>Pros and cons prompts.</li> <li>Role-playing prompts.</li> <li>Loaded/leading questions.</li> <li>Low context prompts.</li> <li>Repeat this.</li> </ul> | <ul style="list-style-type: none"> <li>Context exhaustion: logic-overloading prompts (to reveal unexplainable decisioning processes).</li> </ul> | <ul style="list-style-type: none"> <li>Auto/biographical prompts.</li> <li>Location awareness prompts.</li> <li>Autocompletion prompts.</li> <li>Repeat this.</li> </ul> |
| <b>Manage</b>  |   |   |  |  |



Table G.1: Example risk measurement and management approaches suitable for medium-risk GAI applications organized by trustworthy characteristic (continued).

| Function | Trustworthy Characteristic |                      |                    |
|----------|----------------------------|----------------------|--------------------|
|          | Safe                       | Secure and Resilient | Valid and Reliable |
| Measure  |                            |                      |                    |
| Manage   |                            |                      |                    |

**G.2: Example Medium-risk Generative AI Measurement and Management Plan by Generative AI Risk**

## **Appendix H: Example High-risk Generative AI Measurement and Management Plan**

**H.1: Example High-risk Generative AI Measurement and Management Plan by Trustworthy Characteristic**

**H.2: Example High-risk Generative AI Measurement and Management Plan by Generative AI Risk**